# Equivalence of Single- and Multilocus Markers: Power to Detect Linkage with Composite Markers Derived from Biallelic Loci

Alexander F. Wilson and Alexa J. M. Sorant

Inherited Disease Research Branch, National Human Genome Research Institute, National Institutes of Health, Baltimore

**The reintroduction of biallelic markers, now in the form of single-nucleotide polymorphisms (SNPs), has again raised concerns about the practicality of the use of markers with low heterozygosity for genomic screening for complex traits, even if thousands of such markers are available. Like the early blood-group markers (e.g., Rh and MNS), tightly linked biallelic SNPs can be combined into composite markers with heterozygosity similar to that of short-tandem-repeat polymorphisms. The assumptions that underlie the equivalence between single-locus multiallelic and composite markers are presented. We used computer simulation to determine the power of the Haseman-Elston test for linkage with composite markers when not all of these assumptions hold. The Genometric Analysis Simulation Program was used to simulate continuous and discrete traits, one single-locus four-allele marker, and six biallelic markers. We studied composite markers created from pairs, trios, and quartets of biallelic markers in nuclear families and in independent sib pairs. The power to detect linkage with a two-point approach for composite markers and with a multipoint approach that incorporated all six biallelic markers was compared with that for a single-locus, four-allele reference marker. Although the power to detect linkage with a single biallelic marker was considerably less than that of the reference marker, the power to detect linkage with two- and three-locus composite markers was quite similar to that of the reference marker. The power to detect linkage with four-locus composite markers was similar to that of a multipoint approach.**

## Introduction

The methods used for linkage analysis of discrete and continuous traits in humans have undergone a number of changes since the first traits were linked in humans—the Lutheran and Secretor blood groups (Mohr 1954), the ABO blood group and nail-patella syndrome (Renwick and Lawler 1955), and the Rh blood group and elliptocytosis (Morton 1956). The early markers were primarily blood groups (e.g., the ABO, Rh, MNS, Kidd, Kell, and Lutheran blood groups), and many of these markers were biallelic with low heterozygosity. Some of the most informative early markers, however, were based on systems of tightly linked biallelic loci. The Rh blood-group locus, for example, was specified as three tightly linked biallelic loci: C, D, and E. The MNS blood-group locus was composed of two tightly linked biallelic loci: MN and S. Although the MNS system was usually treated as a single locus without recombination between the component loci, several recombinations between the MN and S loci were reported (Race

and Sanger 1975). The heterozygosity of these systems was often greater than that of single-locus multiallelic markers such as the ABO blood-group locus, depending on the number and frequency of the alleles in the population.

The imminent availability of thousands, or perhaps tens of thousands, of single-nucleotide polymorphisms (SNPs) (Wang et al. 1998) that could be combined into systems of tightly linked biallelic markers (composite markers) offers the promise of providing a large set of highly polymorphic markers that can be used for linkage analysis. Kruglyak (1997) addressed the use of biallelic SNP markers for mapping autosomal dominant disorders with a complete multipoint analysis of all markers on a chromosome; he found that a 1-cM map of biallelic markers (~3,000) provided more information than a 10-cM map of single-locus multiallelic markers—for example, short-tandem-repeat polymorphisms (STRPs) or microsatellites. He noted that "intuitively, one would expect two closely linked biallelics to provide the same information as one microsatellite, and simulations largely confirm this intuition" (Kruglyak 1997, p. 23). However, Hodge et al. (1999) noted that the information available from $k$ SNPs may be less than that from one $2^k$-allele locus and demonstrated this loss of information under the assumption of linkage disequilibrium. In the present study, the assumptions that underlie the equivalence of single-locus multiallelic markers and

composite markers derived from systems of tightly linked biallelic loci were considered for linkage analysis. We used computer simulation to determine the power of tests of linkage with composite markers when not all of these assumptions hold, compared with that for a single-locus multiallelic marker. We also considered the power of a multipoint approach that used information from all six biallelic markers.

## Methods

### Genetic Theory

Consider a single-locus four-allele marker and a composite marker derived from two loci, each having two alleles. Let the alleles for the single-locus marker be denoted $s \in \{1,2,3,4\}$, and let the alleles for the two loci of the composite marker be denoted $h_1 \in \{A,a\}$ and $h_2 \in \{B,b\}$. The four two-locus haplotypes are $h_1 h_2 \in \{AB,Ab,aB,ab\}$. Table 1 illustrates the correspondence between the genotypes for a single-locus four-allele marker and for phase-known and phase-unknown two-locus composite markers. In this instance, there are 10 genotypes for the single-locus and for the phase-known two-locus composite markers, but there are only 9 distinct genotypes for the phase-unknown composite marker. When phase is unknown, the AB/ab and Ab/aB genotypes, which correspond to the 14 and 23 genotypes for the single locus, are indistinguishable. The missing genotype is a result of the inability to distinguish, without additional information, the phase of the double (or more generally, the multiple) heterozygote or heterozygotes. If the alleles of the individual loci are equally frequent, the frequencies of the double heterozygotes (.125) are twice that of the homozygotes (.0625), and only 75% of the population can be assigned an unequivocal genotype.

In this context, single-locus and multilocus marker systems are equivalent under the following assumptions: (1) there is no recombination between the loci of the composite marker; (2) phase can be determined unequivocally; (3) the number of haplotypes at the composite marker corresponds to the number of alleles at the single-locus marker; and (4) the haplotype frequencies for the composite marker correspond to the allele frequencies for the single-locus marker (and the heterozygosities are thus equivalent). For a single-locus four-allele marker, the maximum heterozygosity occurs when the allele frequencies are equal (.25). Similarly, the maximum heterozygosity for each of the two-allele loci occurs when the allele frequencies are equal (.5). If the loci are in linkage equilibrium, the two-locus haplotype frequencies will be .25. If all these assumptions hold, the statistical properties (validity, power, and robustness) of

## Table 1

**Single-Locus Four-Allele Marker Versus Two-Locus Composite Marker**

| Single-Locus Genotype | Two-Locus Genotype | | Frequency[a] |
| | Phase Known | Phase Unknown | |
| --- | --- | --- | --- |
| 11 | AB/AB | AABB | $p_1^2$ |
| 12 | AB/Ab | AABb | $2p_1 p_2$ |
| 13 | AB/aB | AaBB | $2p_1 p_3$ |
| 14 | AB/ab | AaBb[b] | $2p_1 p_4$ |
| 22 | Ab/Ab | AAbb | $p_2^2$ |
| 23 | Ab/aB | AaBb[b] | $2p_2 p_3$ |
| 24 | Ab/ab | Aabb | $2p_2 p_4$ |
| 33 | aB/aB | aaBB | $p_3^2$ |
| 34 | aB/ab | aaBb | $2p_3 p_4$ |
| 44 | ab/ab | aabb | $p_4^2$ |

[a] $p_1$ = frequency of single-locus allele 1 or haplotype AB; $p_2$ = frequency of single-locus allele 2 or haplotype Ab; $p_3$ = frequency of single-locus allele 3 or haplotype aB; $p_4$ = frequency of single-locus allele 4 or haplotype ab.

[b] Indistinguishable phenotypes.

tests of linkage that use single-locus multiallelic and multilocus composite markers are identical.

However, it is not realistic to expect that all of these assumptions hold, particularly if one is screening with biallelic markers in sibship or nuclear-family data in which the phase of the alleles cannot be determined. The extent to which the statistical properties are comparable for the single- and multilocus composite markers depends, in part, on the extent to which the assumptions are true. At best, the composite marker can perform as well as, but no better than, the single-locus marker with the corresponding number of alleles and heterozygosity. At worst, the composite marker can be decomposed into its component loci, and the amount of information will be no worse than that provided by any of the individual single-locus components.

Just as the heterozygosity of a single-locus marker can be increased by increasing the number of alleles (under the assumption that all alleles are equally frequent), the heterozygosity of the composite marker can be increased by increasing the number of loci included in the composite marker. If there are $k$ alleles at a single-locus marker, there are $k(k + 1)/2$ genotypes at that locus. If there are $n$ loci that form an $n$-locus composite marker, then there are $2^n$ haplotypes, $3^n$ $n$-locus genotypes when phase is unknown, and $2^n(2^n + 1)/2$ $n$-locus genotypes when phase is known. The number of genotypes for phase-unknown and phase-known composite markers (and for the corresponding single-locus marker) is presented in table 2. For purposes of comparison, the number of alleles at the single locus is assumed to correspond to the number of haplotypes at the $n$-locus composite

**Table 2**

**Phase-Unknown *n*-Locus Genotypes Compared with Phase-Known *n*-Locus and *k*-Allele Single-Locus Genotypes**

| | | NO. OF *n*-LOCUS GENOTYPES | | |
|---|---|---|---|---|
| NO. OF LOCI | NO. OF ALLELES | Phase Unknown | Phase Known | PHASE RATIO $\kappa$ |
| 1 | 2 | 3 | 3 | 1.00 |
| 2 | 4 | 9 | 10 | .90 |
| 3 | 8 | 27 | 36 | .75 |
| 4 | 16 | 81 | 136 | .60 |
| 5 | 32 | 243 | 528 | .46 |
| 6 | 64 | 729 | 2,080 | .35 |
| 7 | 128 | 2,187 | 8,256 | .26 |
| 8 | 256 | 6,561 | 32,896 | .20 |

marker (i.e., $k = 2^n$). Table 2 also presents the *phase ratio* $\kappa$, defined as the proportion of the number of phase-unknown genotypes to phase-known *n*-locus genotypes: $\kappa = 3^n/[2^n(2^n + 1)/2]$. In the absence of information on phase, the increase in the heterozygosity of the composite marker is not accompanied by equivalent increase in information for linkage analysis. As additional loci are included in the composite marker, the proportion of distinct phase-unknown *n*-locus genotypes relative to phase-known genotypes decreases, as does the proportion of the population that can be assigned an unequivocal genotype. Although the number of *n*-locus genotypes increases exponentially as the number of loci included in the composite marker increases, the increase in the amount of information available for linkage analysis is offset to some extent by the loss of information on phase.

*The Simulation Model*

The Genometric Analysis Simulation Program (G.A.S.P.) version 3.31 (Wilson et al. 1996) was used to simulate a trait locus and seven marker loci: one single-locus four-allele marker ($M_0$) and six biallelic markers ($M_1$, $M_2$, $M_3$, $M_4$, $M_5$, and $M_6$), with map distance and allele frequencies as indicated in table 3. Given the short map interval, the recombination fractions were assumed to be linear, additive, and proportional to the map distance. A continuous variable was used to represent both a continuous trait and an underlying continuous liability for a discrete trait. The trait locus ($T$) was based on a single-locus two-allele additive genetic model with heritabilities (the proportion of the total phenotypic variation due to the single trait locus) .0–.9, with residual variation due to a normally distributed random effect. For the discrete trait, a threshold was set so that the upper 5% of the population would be classified as affected. We performed two sets of simulations. In the first set, 100 nuclear families (each with two parents and four

offspring) were ascertained so that at least two offspring were affected. In this case, the availability of genotyping data on all members of the nuclear family provided partial information on phase. In the second set, 600 independent sib pairs were ascertained so that at least one sib was affected. No parental or other sibling information was available to infer phase. Two thousand replications of each experiment were performed in order to provide a 95% confidence interval of maximum length .02 on the estimate of the type I error rate (.05 ± .01).

*Statistical Genetic Analysis*

A modified version of model-independent two-point sib-pair linkage analysis (Haseman and Elston 1972; Wilson and Elston 1993), as implemented in SIBPAL (S.A.G.E. 1997), was used to test for linkage between the trait $T$ and the single-locus multiallelic marker ($M_0$), each biallelic marker ($M_i$), and all possible pairs ($M_{ij}$), trios ($M_{ijk}$), and quartets ($M_{ijkl}$) of biallelic markers. The squared sib-pair difference was regressed on the estimated proportion of alleles shared identical by descent (IBD), for each marker. Unweighted linear regression was used to test for significance at the .05 level. An affected-sib-pair test was used for the discrete trait, testing whether the proportion of alleles IBD is >.5 for concordantly affected sib pairs only ($P \leq .05$). In a previous series of simulation experiments (data not shown), this affected-sib-pair test had power close to that of the continuous trait in nuclear families of this size, under the same generating model. In the sib-pair tests in SIBPAL, all information in a nuclear family was discarded when inconsistent genotypes were detected (a consequence of a recombination within a composite marker). Although information was removed, the number of degrees of freedom was not correspondingly adjusted downward in the test for the continuous trait, which resulted in a somewhat conservative test.

We also analyzed the continuous trait with a multipoint approach. We used MAPMAKER/SIBS version 2.1 (Kruglyak and Lander 1995) to perform a multipoint version of traditional Haseman-Elston analysis with all six biallelic markers and the true marker map (i.e., the generating model). Because the type I error rate for this

**Table 3**

**Model Parameters for the Simulated Trait Locus (*T*) and Marker Loci (*M_i*)**

| Parameter | $M_0$ | $M_1$ | $M_2$ | $T$ | $M_3$ | $M_4$ | $M_5$ | $M_6$ |
|---|---|---|---|---|---|---|---|---|
| Map distance (cM) | −5 | −3 | −1 | 0 | 1 | 3 | 5 | 7 |
| No. of alleles | 4 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |
| Allele frequencies | .25 | .7 | .8 | .5 | .7 | .8 | .9 | .75 |
| | .25 | .3 | .2 | .5 | .3 | .2 | .1 | .25 |
| | .25 | ... | ... | ... | ... | ... | ... | ... |
| | .25 | ... | ... | ... | ... | ... | ... | ... |

approach was found to be somewhat more liberal than that in SIBPAL, critical values for the $t$ statistics were empirically determined to match the type I error rates obtained for the single-locus marker (.06 and .05 for nuclear families and independent sib pairs, respectively).

In both sets of simulations, 600 sib pairs were used for the sib-pair linkage analysis of the continuous trait. In the first set, each sample comprised 100 families, each with two parents and four offspring (six sib pairs). In the second set, all 600 pairs were independent. Approximately 112 affected sib pairs per sample, on average, were available in the first set for the analysis of the discrete trait. However, considerably fewer independent sib pairs in the second set had both siblings affected; the power of the affected-sib-pair test was based on an average of 28 pairs per sample. Although, in this situation, the absolute power of the test was diminished, the focus was on a comparison of the power of the tests of multilocus composite markers relative to that of the single-locus multiallelic marker, not on the absolute power of the test.

The proportion of samples that gave an indication of linkage was determined for each marker, for trait heritabilities of .0–.9 in increments of .1. When the heritability was 0—that is, when there was no genetic component underlying the trait—the proportion of samples that gave an indication of linkage can be taken as an estimate of the type I error rate—that is, the proportion of false positives. For all other heritabilities, this proportion estimates the power of the test to detect linkage under the specific generating model.

## Results

The power to detect linkage with a two-point approach, relative to that of a single-locus four-allele marker, was averaged over nine heritabilities (.1–.9) and is presented, in table 4, as an average and range over individual biallelic markers ($M_i$) and over sets of two-, three-, and four-

locus composite markers ($M_{ij}$, $M_{ijk}$, and $M_{ijkl}$, respectively). Similarly, the power to detect linkage with a multipoint approach ($M_{multipoint}$) was compared with the power of the single-locus marker. The absolute power to detect linkage for the single-locus four-allele marker ($M_0$) is given as a reference. We present relative power for both continuous and discrete traits, for samples with phase partially known (nuclear families) and with phase unknown (independent pairs).

For the continuous trait with nuclear-family data, there was, on average, a 16% relative increase in power when we used a two-locus composite marker, compared with that of a one-locus biallelic marker, a 4% relative increase in power when we used a three-locus composite marker, compared with that of a two-locus composite marker, and a 2% relative increase in power when we used a four-locus composite marker, compared with that of a three-locus composite marker. Similarly, for the discrete trait, the increases were 30%, 8%, and 3%, respectively; with independent sib pairs, comparable increases were 14%, 5%, and 2%, respectively, for the continuous trait and 18%, 8%, and 3%, respectively, for the discrete trait.

For nuclear-family data, the type I error rate was .04–.07, averaging .058 and .055 for continuous and discrete traits, respectively. For independent sib-pair data, the rate was .02–.09, averaging .051 and .045 for continuous and discrete traits, respectively.

## Discussion

When taken individually, the power to detect linkage with biallelic markers with heterozygosity typical of SNPs was less than that detected with a single-locus multiallelic marker—sometimes substantially so. This was true even when the individual biallelic markers were much closer to the trait locus (1–3 cM) than was the single-locus multiallelic marker (5 cM). In nuclear-family data, power was as much as 34% and 52% less than

**Table 4**

**Average Power of Biallelic Markers (Single and Composite) and Multipoint Tests, Relative to the Power of a Single-Locus Four-Allele Marker**

| MARKER | NUCLEAR FAMILIES | | INDEPENDENT PAIRS | |
|---|---|---|---|---|
| | Continuous Trait | Qualitative Trait | Continuous Trait | Qualitative Trait |
| | Absolute Power of the Single-Locus Four-Allele Marker | | | |
| $M_0$ | .84 | .78 | .84 | .36 |
| | Average and Range of Power Relative to $M_0$ | | | |
| $M_i$ | .83 (.66–.92) | .72 (.48–.86) | .85 (.72–.92) | .71 (.54–.85) |
| $M_{ij}$ | .96 (.89–1.01) | .93 (.79–1.02) | .96 (.90–1.00) | .84 (.69–1.00) |
| $M_{ijk}$ | 1.00 (.97–1.03) | 1.01 (.95–1.07) | 1.01 (.98–1.03) | .90 (.75–1.06) |
| $M_{ijkl}$ | 1.02 (1.00–1.04) | 1.04 (1.00–1.07) | 1.03 (1.02–1.05) | .93 (.82–1.06) |
| $M_{multipoint}$ | 1.02 (1.00–1.03) | ... | 1.03 (1.02–1.03) | ... |

that of the single-locus multiallelic marker for continuous and discrete traits, respectively, and occurred at the biallelic marker with the lowest heterozygosity ($M_5$), at a distance of 5 cM from the trait locus. For independent sib pairs, the power was 28% and 46% less for the continuous and discrete traits, respectively, at that same marker. The power obtained with biallelic markers in this study was consistent with determinations of power from a number of early simulation studies of model-independent sib-pair linkage analysis (Blackwelder and Elston 1982; Amos et al. 1989).

When composite markers were formed from pairs or trios of nearby biallelic markers (with recombination within the composite marker ignored), the power to detect linkage was quite similar to that obtained with a single-locus four-allele marker, for both continuous and discrete traits, when nuclear-family data were considered, and for the continuous trait, when independent sib pairs were used. The addition of a fourth biallelic locus to the composite marker was necessary to bring the power close to that of the four-allele single-locus marker, for the discrete trait for independent sib pairs, although this may be because there were substantially fewer pairs used in the analysis, owing to the method of ascertainment.

Although, in general, the relative power increased as additional loci were added to the composite marker, the improvement in power decreased as the number of loci included in the composite marker increased. The diminishing improvement may be due to an increased probability of recombination within the composite marker and to the decrease in the number of distinct phase-unknown $n$-locus genotypes that could be identified, relative to the number of distinct phase-known $n$-locus genotypes (i.e., decreasing $\kappa$). Holmans and Clayton (1995) noted a similar effect in a study of the efficiency of linkage when they used an affected-sib-pair approach. They attributed this effect to haplotypes with phase uncertainties.

Although the power to detect linkage with three- and four-locus composite markers was compared with that of a single-locus four-allele locus to illustrate the diminishing improvement in power, a more appropriate comparison would have been between the three- and four-locus composite markers and single loci with 8 and 16 alleles, respectively. Given the modest improvement over a single-locus four-allele marker, a substantial loss of power with three- and four-locus composites would be expected, compared with that for single loci with the corresponding number of alleles. However, if there is no multiallelic single-locus marker with sufficient heterozygosity in the interval to be searched, a three- or four-locus composite of biallelic markers would be useful, particularly if information on phase is known.

It is important to note that, in linkage analysis, it is the heterozygosity of the marker or of the region (regional heterozygosity) that determines the amount of information available, not just the density of the marker map. In the limited situations considered, the power to detect linkage with a two-point approach for four-locus composite markers (over a 6–10-cM map) compares favorably with that of a multipoint approach that uses information from all six biallelic markers. Although a likelihood-based multipoint approach that incorporates all the marker information over large chromosomal regions would be the most powerful approach, the amount of computation increases disproportionately as the size of the families studied increases, which makes multipoint analysis of very large families problematic. Furthermore, a standard multipoint approach assumes that the order of the markers and the distances between them are known with certainty. Thus, until the SNP map is better characterized, this composite-marker two-point approach may be useful as an approximation to standard multipoint methods. This approach should be more robust with respect to errors in map order and distance and should be computationally trivial, compared with standard multipoint methods. This approach could also be combined with a moving-average approach to significance testing (Goldin and Chase 1997; Goldin et al. 1999), which may prove to be an effective alternative over the near term.

There was a slightly higher than expected false-positive rate in these tests, although the increase was within the 95% confidence interval. The Haseman-Elston test appeared to be statistically valid even when there was recombination within the composite marker. The inflated type I error rate was most pronounced for continuous traits in nuclear families and was probably due to nonnormality and reduction in the variation of the trait because the families were ascertained for sibs with extreme values.

There are, however, some aspects of the combination of biallelic loci into composite markers that are problematic. Mistyping rates will be compounded if single loci are combined. A mistyping rate of $r$ per single locus corresponds to a mistyping rate at the composite marker of $1 - (1 - r)^n$, where $n$ is the number of loci in the composite marker. If the mistyping rate for a single locus was 0.5%, for example, the mistyping rates for two-, three-, and four-locus composite markers would be 1%, 1.5%, and 2%, respectively. If the mistyping rate for a single locus was as much as 5%, however, the corresponding mistyping rates would be 9.8%, 14.3%, and 18.5%, respectively. A similar problem occurs with missing data, although missing genotypes or portions of genotypes can sometimes be inferred.

On the basis of this study, it appears that composite markers based on two to three loci with heterozygosity typical of SNPs (~.3) and uniformly distributed on a 2-

cM map appear to be nearly as polymorphic as STRPs and provide almost the same amount of information for model-independent linkage analysis, even in the absence of information on phase. A screening strategy that uses a 2-cM SNP map (~1,600–1,700 biallelic loci) should be sufficient to provide nearly the same information for linkage analysis as is provided by current genomic screens that use 350–400 STRPs, provided that the mistyping rates for the composite markers are comparable to those for single-locus STRPs. Although retyping of samples for inconsistent or failed results has, for decades, been common practice with gel-based electrophoretic technologies, retyping with a chip-based technology is more problematic and may require either the use of a second chip or of a set of duplicate markers on each chip.

## Acknowledgments

## References

Amos CI, Elston RC, Wilson AF, Bailey-Wilson JE (1989) A more powerful robust sib-pair test of linkage for quantitative traits. Genet Epidemiol 6:435–449

Blackwelder WC, Elston RC (1982) Power and robustness of sib-pair linkage tests and extension to larger sibships. Commun Stat Theor Methods 11:449–484

Goldin LR, Chase GA (1997) Improvement of the power to detect complex disease genes by regional inference procedures. Genet Epidemiol 14:785–789

Goldin LR, Chase GA, Wilson AF (1999) Regional inference with averaged p values increases the power to detect linkage. Genet Epidemiol 17:157–164

Haseman JK, Elston RC (1972) The investigation of linkage between a quantitative trait and a marker locus. Behav Genet 2:3–19

Hodge SE, Boehnke M, Spence MA (1999) Loss of information due to ambiguous haplotyping of SNPs. Nat Genet 21:360–361

Holmans P, Clayton D (1995) Efficiency of typing unaffected relatives in an affected-sib-pair linkage study with single-locus and multiple tightly linked markers. Am J Hum Genet 57:1221–1232

Kruglyak L (1997) The use of a genetic map of biallelic markers in linkage studies. Nat Genet 17:21–24

Kruglyak L, Lander E (1995) Complete multipoint sib-pair analysis of qualitative and quantitative traits. Am J Hum Genet 57:439–454

Mohr J (1954) A study of linkage in man. Munksgaard, Copenhagen

Morton NE (1956) The detection and estimation of linkage between the genes for elliptocytosis and the Rh blood type. Am J Hum Genet 8:80–96

Race RR, Sanger R (1975) Blood groups in man. Oxford University Press, Oxford

Renwick JH, Lawler SD (1955) Genetical linkage between the ABO and nail-patella loci. Ann Hum Genet 19:312–331

S.A.G.E. (1997) Statistical analysis for genetic epidemiology, release 3.1. Computer program package available from the Department of Epidemiology and Biostatistics, Rammelkamp Center for Education and Research, MetroHealth Campus, Case Western Reserve University, Cleveland

Wang DG, Fan J-B, Siao C-J, Berno A, Young P, Sapolsky R, Ghandour G, et al (1998) Large-scale identification, mapping and genotyping of single-nucleotide polymorphisms in the human genome. Science 280:1077–1082

Wilson AF, Bailey-Wilson JE, Pugh EW, Sorant AJM (1996) The Genometric Analysis Simulation Program (G.A.S.P.): a software tool for testing and investigating methods in statistical genetics. Am J Hum Genet Suppl 59:A193

Wilson AF, Elston RC (1993) Statistical validity of the Haseman-Elston sib-pair test in small samples. Genet Epidemiol 10:593–598